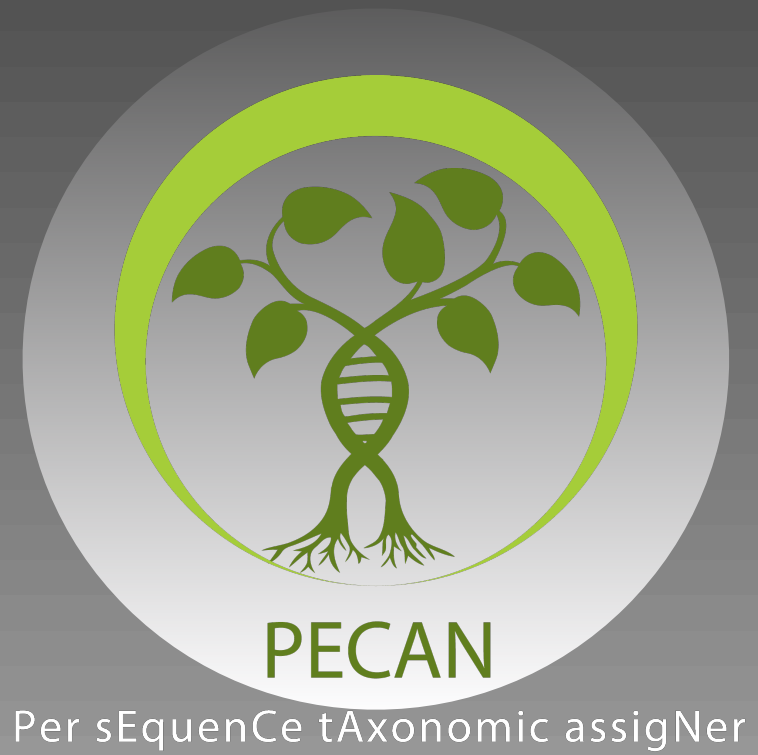


PECAN

A fast, novel 16S rRNA gene sequence non-clustering based taxonomic assignment tool

Johanna B. Holm, Pawel Gajer, Jacques Ravel
Institute for Genome Sciences, University of Maryland Baltimore School of Medicine



Background

Clustering of sequences into Operational Taxonomic Units (OTUs) has become a mainstream approach to facilitate taxonomic classification of large numbers of 16S rRNA gene sequences. This is partly due to the high computational requirements for processing each sequence in increasingly large datasets. A primary focus of the field has been development and improvement of OTU-based sequence clustering methods that rely on distances between each pair of sequences in a dataset. Following OTU-based clustering, representative sequences are commonly classified using tools such as the RDP Naïve Bayesian Classifier (Wang *et al.* 2007), and the resulting classification transitively assigned to all sequences comprising that OTU. However, problems with this strategy exist (Nguyen *et al.*, 2016). Here, we present **PECAN**, a novel per sequence taxonomic assigner which quickly and accurately classifies millions of 16S rRNA gene sequences using higher order Markov Chain models built from a user-specified set of reference sequences.

Methods & Usage

1. Obtain, align, & truncate reference sequences to variable region(s) of choice

- RDP database, MAFFT v 7.222, & mothur v 1.36.0

2. Produce phylogenetic tree

- FastTree JC + CAT

3. Curate taxonomic annotation curation (Figure 1)

- An internal pipeline to automatically assign corrected taxonomic annotations.

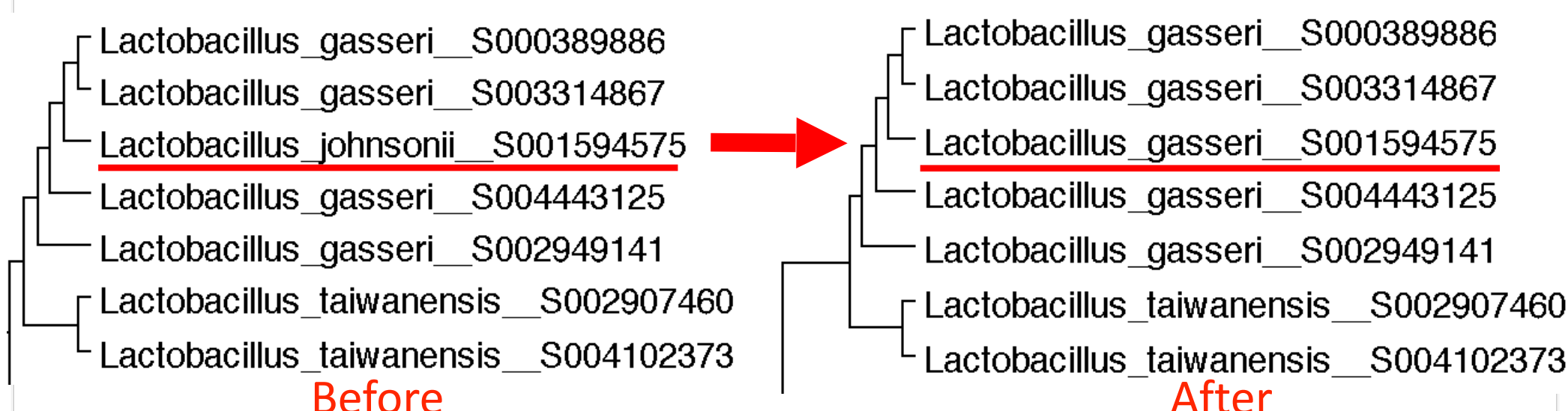


Figure 1. Example of taxonomic annotation curation. Prior to curation, sequence S001594575 was annotated as *L. johnsonii*. It is corrected to *L. gasseri* after alignment-dependent curation.

4. Curation produces a taxonomic tree containing all sequences (Figure 2). For each species, genus, etc., construct 7th order Markov Chain model (Table 1).

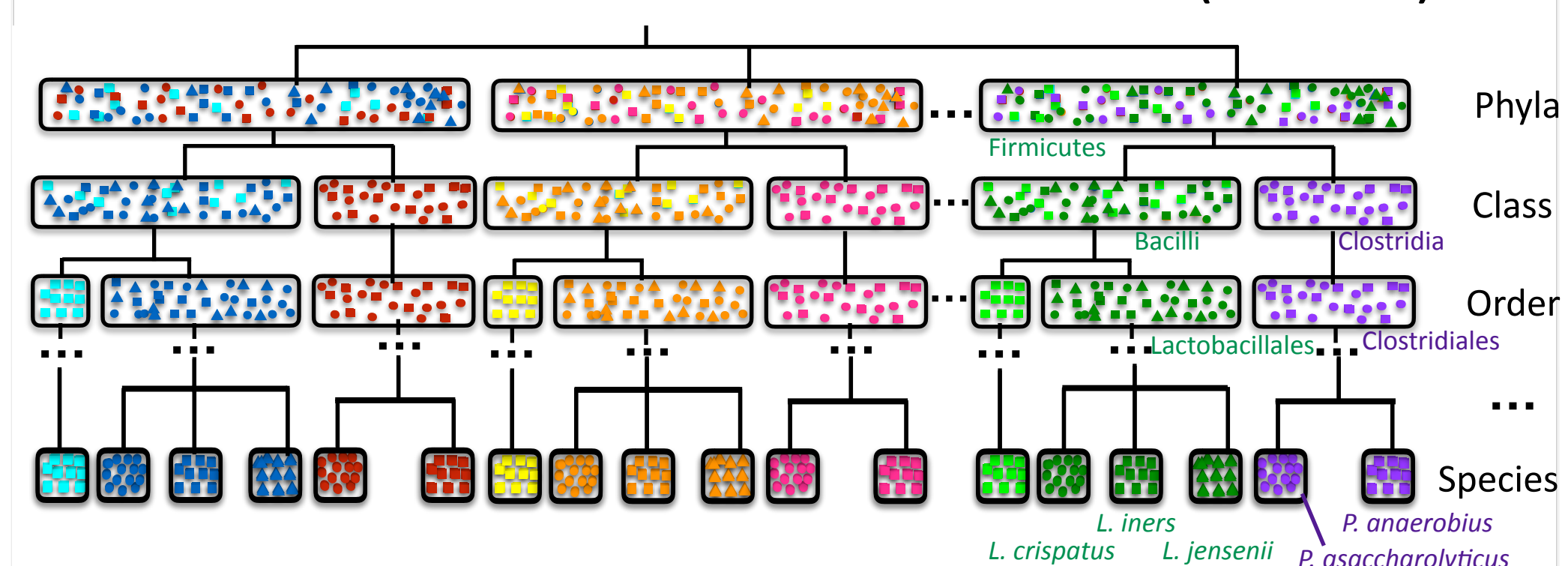


Figure 2. A taxonomic tree is produced from the alignment of reference sequences

Table 1. Markov Chain models of the probability of observing a particular base following a particular 7mer using *Lactobacillus* spp. as an example.

	A1 7mer _a	C1 7mer _a	G1 7mer _a	T1 7mer _a	A1 7mer _b	C1 7mer _b
<i>L. iners</i>	P(A 7mer _a)	P(C 7mer _a)	P(G 7mer _a)	P(T 7mer _a)	P(A 7mer _b)	P(C 7mer _b)
<i>L. jensenii</i>	P(A 7mer _a)	P(C 7mer _a)	P(G 7mer _a)	P(T 7mer _a)	P(A 7mer _b)	P(C 7mer _b)
<i>L. crispatus</i>	P(A 7mer _a)	P(C 7mer _a)	P(G 7mer _a)	P(T 7mer _a)	P(A 7mer _b)	P(C 7mer _b)
...	P(A 7mer _a)	P(C 7mer _a)	P(G 7mer _a)	P(T 7mer _a)	P(A 7mer _b)	P(C 7mer _b)
<i>g_Lactobacillus</i>	P(A 7mer _a)	P(C 7mer _a)	P(G 7mer _a)	P(T 7mer _a)	P(A 7mer _b)	P(C 7mer _b)
...
<i>f_Lactobacillaceae</i>	P(A 7mer _a)	P(C 7mer _a)	P(G 7mer _a)	P(T 7mer _a)	P(A 7mer _b)	P(C 7mer _b)
...
<i>o_Lactobacillales</i>	P(A 7mer _a)	P(C 7mer _a)	P(G 7mer _a)	P(T 7mer _a)	P(A 7mer _b)	P(C 7mer _b)
...
<i>c_Bacilli</i>	P(A 7mer _a)	P(C 7mer _a)	P(G 7mer _a)	P(T 7mer _a)	P(A 7mer _b)	P(C 7mer _b)
...
<i>p_Firmicutes</i>	P(A 7mer _a)	P(C 7mer _a)	P(G 7mer _a)	P(T 7mer _a)	P(A 7mer _b)	P(C 7mer _b)

"For *L. iners*, this is the probability of observing a C given the previous 7 nucleotides" (here, labeled as "7mer_a")

5. The reference taxonomic tree is represented by Markov Chain models (Figure 3).

